

GHAJAR EXHIBIT 37

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 1

IN THE UNITED STATES DISTRICT COURT
FOR THE NORTHERN DISTRICT OF CALIFORNIA

_____ RICHARD KADRY, et al,)	CASE NO.
)	3:23-cv-03417-VC
Plaintiff,)	
-against-)	
META PLATFORMS, INC.)	
Defendant.)	
_____)	

HIGHLY CONFIDENTIAL

UNDER THE PROTECTIVE ORDER

VIDEO-RECORDED DEPOSITION OF
MICHAEL CLARK 30 (B) (6)
Zoom Video Conferencing
12/19/2024
3:06 p.m. (MST)

REPORTED BY: MONIQUE CABRERA

DIGITAL EVIDENCE GROUP
1730 M. Street, NW, Suite 812
Washington, D.C. 20036
(202) 232-0646

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 25

1 way more effective to put other mitigations in
2 place earlier in the process from a pretraining
3 perspective in how the model was trained and then
4 also to measure from an ongoing perspective
5 memorization and the likelihood that content
6 could be repeated from training data versus
7 trying to build an exact match filter list like
8 this.

9 Q. Thank you, Mr. Clark.

10 I see a list of eight rows of eight
11 types of mitigations on this document. You have
12 called this document on several occasions an
13 exploratory document. So I would like to have
14 Meta tell me, which of these eight mitigations
15 were not implemented?

16 A. Okay. Let me read through them
17 again real quick for that. We can go through
18 them one by one.

19 Number one, deduplication of
20 datasets was adopted as a mitigation. And
21 there's a mitigation as a requirement from the
22 dataset review set process and before data could

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 26

1 be used for training.

2 The statement "Single pass: We
3 train with one epoch over IP rich data" -- and
4 then in parentheses (common crawl) which has
5 shown in academic literature to reduce
6 memorization," that was adopted as a mitigation.

7 For clarity, because we were also
8 measuring the likelihood of memorization doing
9 just a single pass or single epoch either,
10 that -- that expanded over time because we were
11 able to demonstrate that it was lower numbers of
12 epochs, not just a single pass. So it's slightly
13 modified from how this is written.

14 The basic decoding methods where "We
15 applied decoding methods (sampling, temperature)
16 which statistically reduce likelihood of
17 memorization," those are applied naturally in
18 the -- or in the response side of the model.
19 They are not identified as specific mitigations.

20 And why you would set the decoding
21 methods had more to do with prompt quality
22 response or -- or response quality as opposed to

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 27

1 IP mitigations so it's not specifically
2 identified as an IP mitigation, but we do have --
3 we do have optimized decoding methods that are
4 not at zero.

5 Sampling and temperature have to do
6 with the randomization of responses, and so we --
7 we don't want them fully random, but we don't
8 want -- we don't want them under-randomized.

9 Number 4 is purchased data. We have
10 added our own purchased data, human annotations,
11 license data into the training mix. This -- from
12 a IP mitigation perspective, this was not adopted
13 as an IP mitigation.

14 We do have our own annotation data
15 that we have added to models and depend on that
16 for a variety of portions of the training data.

17 License data, we don't have any
18 textual licensed data or purchased licensed data.

19 Targeted URL removal where it says
20 "We can delete business risk data sources by URL
21 from our pretraining date (i.e., ticketmaster,
22 orbitz or booking.com), this one, the way it's

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 28

1 written is fairly poorly written.

2 What this is talking about is from
3 datasets that are based on URLs, which a good
4 example would be Common Crawl where Common Crawl
5 stores the data by URL, what data came from
6 Ticketmaster or what data came from Orbitz or
7 booking.com.

8 And so in the processing of that
9 training data, it's possible to say there are
10 certain domains which we would not want to
11 process. And so by targeted URL removal or the
12 block list that is referred to the risks and
13 blockers -- we had talked about the block list, I
14 believe prior in my testimony. What this blocked
15 list referred to were like Ticketmaster, Orbitz
16 and booking.com specifically are all sites who --
17 who had strong enforcement around scraping of
18 their sites.

19 And this is, once again being in
20 March of '23, prior to Common Crawl being very
21 transparent about what they crawled and didn't
22 and how people could block Common Crawl from

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 29

1 crawling their data.

2 And so we filtered out sites that
3 had -- when we processed the Common Crawl data,
4 we filtered out the data from sites whose URLs
5 matched against that block list because of how
6 the data was scraped.

7 Number 6 is the hotfixed IP filter
8 list.

9 So just to clarify on 5, but 5 was
10 adopted but not in the form in which is it's
11 currently written or labeled or titled.

12 6, the hotfixed IP filter list was
13 not adopted as a mitigation.

14 7, "Generation time duplicate
15 detection: build an n-gram index of the training
16 set to search for n-gram overlap at generation
17 time." And "generations over a particular
18 threshold could be rejected."

19 This was not adopted as a
20 mitigation, both from a technical feasibility
21 perspective but also because we put investments
22 into sampling of the data instead and measuring

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 30

1 memorization potential of the model end to end.

2 And then finally, "Dynamic sampling:
3 based on model likelihood, reject select
4 generations with the hope that" those "are more
5 likely to match training corpus."

6 Once again, rejected because it's --
7 it's not clear that it would have made any
8 difference, and we focused on the mitigations
9 both from a pretraining perspective in addition
10 to measuring the likelihood of -- of memorization
11 from an output or from a reidentified --
12 reidentified ability perspective.

13 Q. Thank you, Mr. Clark.

14 What copyright/IP mitigations were
15 adopted that are not on this list?

16 A. We had deduplication. Sorry. I
17 have got to walk through them just to remember
18 them all.

19 Single pass, we had measurement of
20 memorization. We had content from the block
21 list, being URLs that matched the block list for
22 datasets who are based around specific URLs.

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 31

1 And I just need a moment to
2 recollect.

3 Oh, and data that self-identified as
4 pirated or stolen was to be removed or excluded
5 from training. A lot of these -- actually all of
6 these mitigations, except for the last one,
7 originated from doing this for privacy and
8 preventing the model from regurgitating or
9 reidentifying information about private
10 individuals.

11 And so these overlapped with that,
12 and then the measurement and memorization being
13 the final one.

14 Q. And in terms of the memorization
15 issue, you mentioned that, you know, it's an
16 exploratory document. And Kathleen can correct
17 me if I am mischaracterizing your testimony, but
18 I think your testimony was that Meta's current
19 position is that -- is not that you can consider
20 any memorization to be a copyright risk.

21 Is there an amount of memorization
22 that Meta does consider to be a copyright risk?

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 32

1 MS. HARNETT: Object to form.

2 A. Once again, "memorization" is a very
3 generic term. Some memorization is intentional,
4 and so some data is trained to have a higher
5 likelihood of replying the information. Is
6 Abraham Lincoln a president? Do you expect the
7 model to be able to know that?

8 The threshold that was put in place
9 is that Llama 2 with low levels of memorization,
10 that for future models or future training,
11 that -- that future models could not have a
12 higher level of memorization than what Llama 2
13 was. And so Llama 2 became the -- the lowest
14 threshold for potential memorization.

15 Q. And here in this document, there is
16 actually a definition of memorization. Could you
17 look to that with the three bullets on the first
18 page ending in 459?

19 A. Yes. And so this is under the
20 definition heading that states, "On that note, we
21 define memorization as" -- and then a first
22 bullet -- "An uncommon phrase (as determined by

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 33

1 perplexity)."

2 Second bullet, "At least 50 words in
3 length." And the third bullet being "With a 90%"
4 match to the training data.

5 Q. Does Meta agree that that's the
6 working definition of memorization?

7 A. No. Once again, the purpose of this
8 document was a theoretical exploration of
9 potential mitigations and potential definitions
10 by a product manager and engineers that then fed
11 into what -- what memorization could be.

12 This, once again, is kind of, like,
13 I'm defining any memorization to be a copyright
14 risk and memorization is a very generic term.
15 And so this -- this does not what was -- this is
16 not what was Meta's position, as far as
17 memorization.

18 Q. So what is Meta's definition of
19 memorization if it's not this?

20 MS. HARNETT: I object to the --
21 object to the form and the scope and also,
22 foundation.

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 34

1 A. Memorization is -- memorization, as
2 defined, is not these three items where it's an
3 uncommon phrase or at least 50 words in length or
4 with a 90 percent match to training data.

5 As it was defined, memorization is,
6 once again, a broad technological concept. And
7 so memorization is defined as the model's ability
8 to repeat what is in training data as an exact
9 match and is measured across various lengths --
10 not 50 words in length, but it's measured across
11 a variety of the different lengths of phrases and
12 random sets in order to do that testing.

13 And memorization -- memorization is
14 generically defined as that ability and then --
15 yeah, that's -- there isn't a more specific
16 definition.

17 BY MR. STEIN:

18 Q. I just don't understand the problem
19 with this definition that's in this document.

20 What is Meta saying is incorrect
21 about this definition?

22 MS. HARNETT: Objection to the form.

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 35

1 A. Sure. An uncommon phrase isn't
2 measurable. It's not determined -- it has
3 nothing to do with -- there are many uncommon
4 phrases in training data. And the model in
5 generating out -- outputs, both for how it's
6 instructed to do so from a certain amount of
7 randomness, but also how it is trained to do so
8 can generate an uncommon phrase.

9 And so just the way this is worded
10 is -- doesn't have anything to do with the
11 likelihood of memorization or being memorization.

12 At least 50 words in length, as I
13 stated, when we look for -- we have defined
14 memorization to be the likelihood of the model
15 outputting content that matches what is in
16 training data. We don't just look at a specific
17 length of data. We measure memorization across
18 shorter and longer sets of phrases, numbers of
19 words.

20 And then with the 90 percent match
21 to the training data, once again, because the
22 other two bullets here are either inaccurate or

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 36

1 ineffective or only, like, points in time in
2 estimation for this explorative study, the -- the
3 90 percent match to training data doesn't match
4 the model's ability to output training data
5 exactly.

6 So these -- these three criteria are
7 not part of what the application is in Meta's
8 position, which is why I state, like,
9 memorization, as defined, is the model's ability
10 to output data exactly as it was in a model.

11 BY MR. STEIN:

12 Q. So I -- if I ask Llama to tell me
13 the first paragraph of Stephen King's book The
14 Stand, it will tell me what that first paragraph
15 is. Is that memorization?

16 MS. HARNETT: Objection to the form
17 and the scope. And this is not about the
18 document.

19 MR. STEIN: Well, we're trying to
20 understand the definition of memorization,
21 which is the document. And he's telling me
22 that this definition isn't -- isn't Meta's

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 43

1 around broad, publicly availability data, we --
2 you know, it's -- nothing is 100 percent, but we
3 put the -- in place the mitigations to prevent as
4 much memorization as possible.

5 BY MR. STEIN:

6 Q. Does Meta consider a removal of
7 copyright information to a part of copyright IP
8 mitigations?

9 MS. HARNETT: Object to the form.
10 Vague and not related to the document.

11 MR. STEIN: Is there -- just --

12 MS. HARNETT: Can you tie that to
13 the document?

14 BY MR. STEIN:

15 Q. The document is about copyright IP
16 mitigations. So I'm wondering if Meta would
17 consider the removal of copyright information
18 or -- or copyright management information to be a
19 copyright IP mitigation.

20 MS. HARNETT: Object to the form.
21 Vague and not related to the document.

22 He is not here to testify about

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 44

1 mitigations generally. He is here to testify
2 about this document.

3 MR. STEIN: Are you instructing him
4 to not answer this particular question?

5 MS. HARNETT: He can answer the
6 question because we want to get the
7 deposition over, but I am just -- I am making
8 my objection for the record that you are not
9 asking questions about the document or
10 directly related to the document.

11 A. We had discussed this prior in a
12 prior deposition, and I did not prepare or have
13 that document as a reference for all of the types
14 of information. Copyright is as a phrase or the
15 copyright symbol or the back/N or the similar
16 things are repeatable sets of ASCII characters
17 that were not removed because they were copyright
18 material. They were removed because of the
19 amount of times that that and other kinds of
20 things were removed from the content because of
21 how repeatable they were, how they existed in
22 spots that would interrupt tokens and break

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 45

1 tokens.

2 And without doing that cleanup from
3 a data parsing perspective, that would equal poor
4 performance in the model and not allow the model
5 to hit the industry benchmark. So extra line
6 spaces, for instance, where somebody hits return
7 five times, those five were removed so the
8 paragraph ending is still there, but the extra
9 things are not.

10 So it removes noise and it removes
11 repeatability and it removes other things.
12 Otherwise, the model would just regurgitate the
13 things that were repeated the most.

14 Q. And in terms of copyright/IP
15 mitigations, when Meta measures memorization,
16 does it have to call back on the training
17 datasets themselves?

18 MS. HARNETT: Object to form and
19 beyond the scope and vague.

20 A. The -- as we walked through in prior
21 testimony, the memorization process pulls random
22 strings of ASCII characters in various lengths

12/19/2024

Richard Kadrey, et al. v. Meta Platforms, Inc. Michael Clark 30(b)(6)
Highly Confidential - Under the Protective Order

Page 58

1 CERTIFICATE OF SHORTHAND REPORTER

2 NOTARY PUBLIC

3 I, Monique Cabrera, the officer
4 before whom the foregoing deposition was
5 taken, do hereby certify that the foregoing
6 transcript is a true and correct record of
7 the testimony given; that said testimony was
8 taken by me stenographically and thereafter
9 reduced to typewriting under my direction;
10 and that I am neither counsel for, related
11 to, nor employed by any of the parties to
12 this case and have no interest, financial or
13 otherwise, in its outcome.

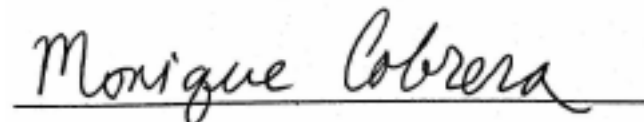
14 IN WITNESS WHEREOF, I have hereunto
15 set my hand this 19th day of December, 2024.

16

17

18

19



MONIQUE CABRERA

20 Notary Public in and for the State of New York
County of Suffolk

21 My Commission No. 01CA6043156

22 Expires: 06/12/2026